

1 Abstract

We present **dSPRINT**: domain Sequence-based PRediction of INTeraction sites, an ensemble of machine learning classifiers using a novel stacking architecture, that predict binding positions within protein domains.

2 Background

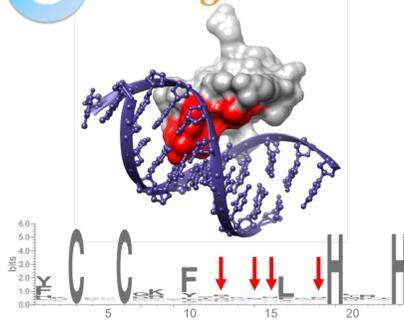


Fig 1: SDPs exemplified in the Cys₂-His₂ Zinc finger domain. Known¹ binding positions are colored in red on the domain surface², and pointed by red arrows on the domain sequence logo³. They are critical for the DNA-binding specificity, and are not conserved.

4 Results: Global evaluation

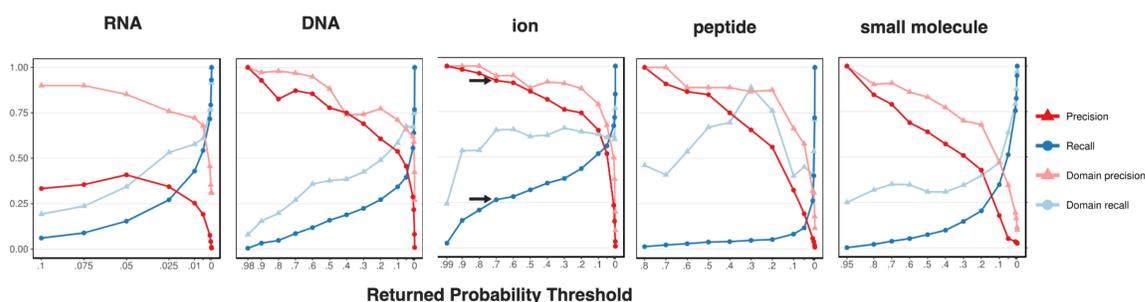


Fig 3: Prediction scores precision and recall curves. Supervised learning using InteracDome⁴ structure-based positional scores, in 5 folds cross-validation.

	Number of most reliable predictions chosen per domain	Number of domains with at least one Correct prediction	Fraction of domains With correct predictions
RNA	1	13	62%
	3	16	76%
	5	18	86%
DNA	1	22	67%
	3	27	82%
	5	29	88%
ion	1	47	52%
	3	60	66%
	5	64	70%
peptide	1	29	40%
	3	44	61%
	5	48	67%
small molecule	1	72	55%
	3	96	73%
	5	104	79%

Table 1: The most reliable prediction(s) in each binding domain.

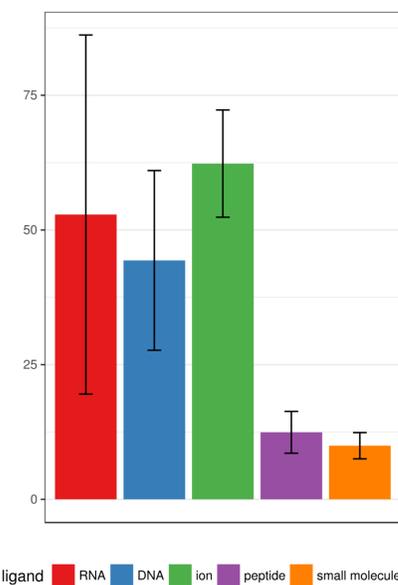


Fig 4: Ligand AUPRC fold improvement. The ratio of the AUPRC to an average baseline corresponding to the fraction of binding and neutral positions of that ligand at that CV fold.

5 Results: Per-domain evaluation



Fig 5: Performance evaluation on ligand-binding domains. The table represents domain-ligand pairs with performance exceeding that of the random baselines (the dashed lines) of AUC=0.5 and AUPRC fold ratio=1.

3 Methods: ML stacked architecture

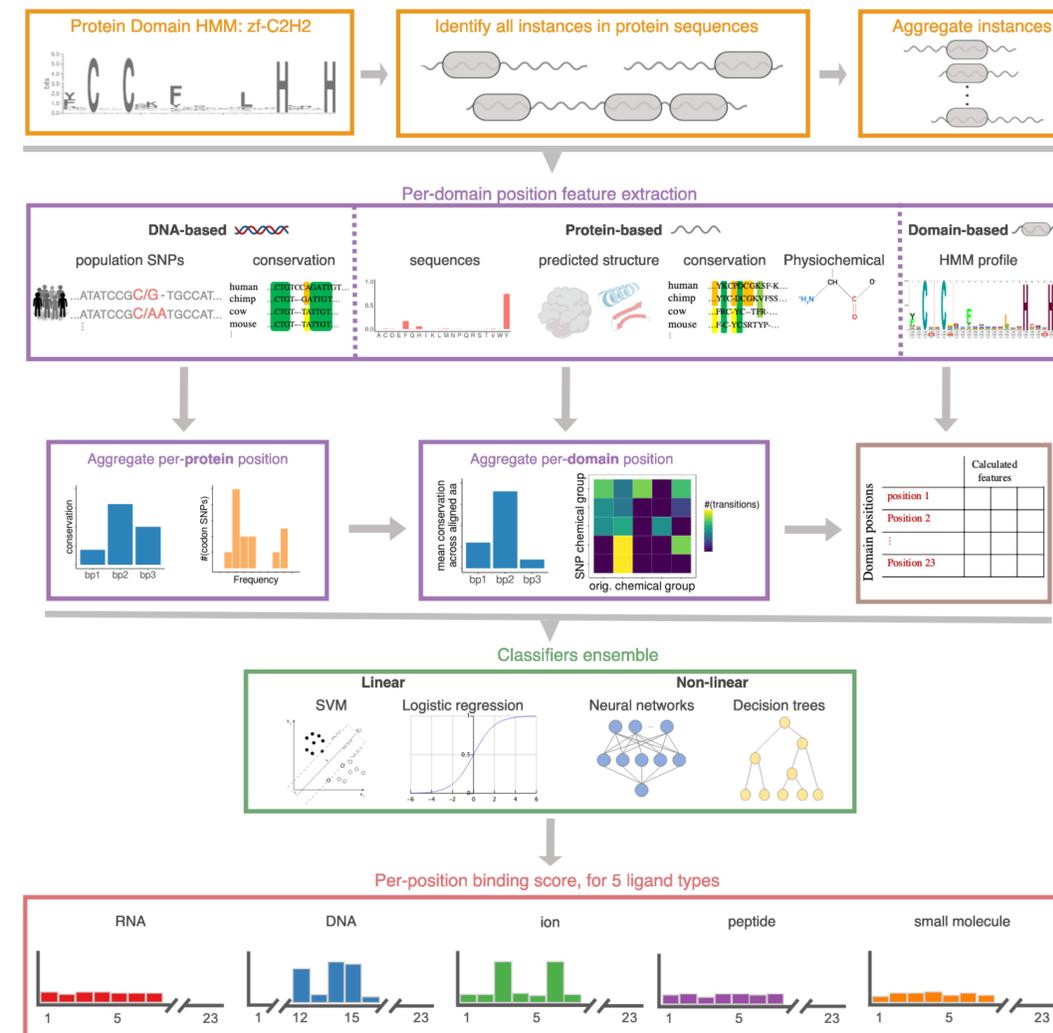


Fig 2: dSPRINT workflow for domain-centered per-position prediction.

6 Significance

Systematic identification of residues essential for ligand-binding would have a farther-reaching applications:

- Identify the functional impact of coding variants
- Explore the variation that evolves in protein interaction network
- Characterize the effect of mutations in the context of disease
- Suggest molecular targets for therapeutic intervention

References

- [1] Wolfe et al. Ann. Rev. Biochem. (2001)
 [2] Pettersen et al. J. Comput. Chem. (2004)
 [3] Crooks et al. Genome Res. (2004)
 [4] Kobren and Singh, NAR (2018)

Acknowledgements

We would like to thank Shilpa N. Kobren, and the Singh lab for their helpful insights. Funded by NIH GM076275 (MS), T32 HG003284 (AE), and CA208148 (MS).